

P&S Processing-in-Memory

Real-World Processing-in-Memory Architectures:
SK Hynix Accelerator-in-Memory

Dr. Juan Gómez Luna

Prof. Onur Mutlu

ETH Zürich

Fall 2022

15 November 2022

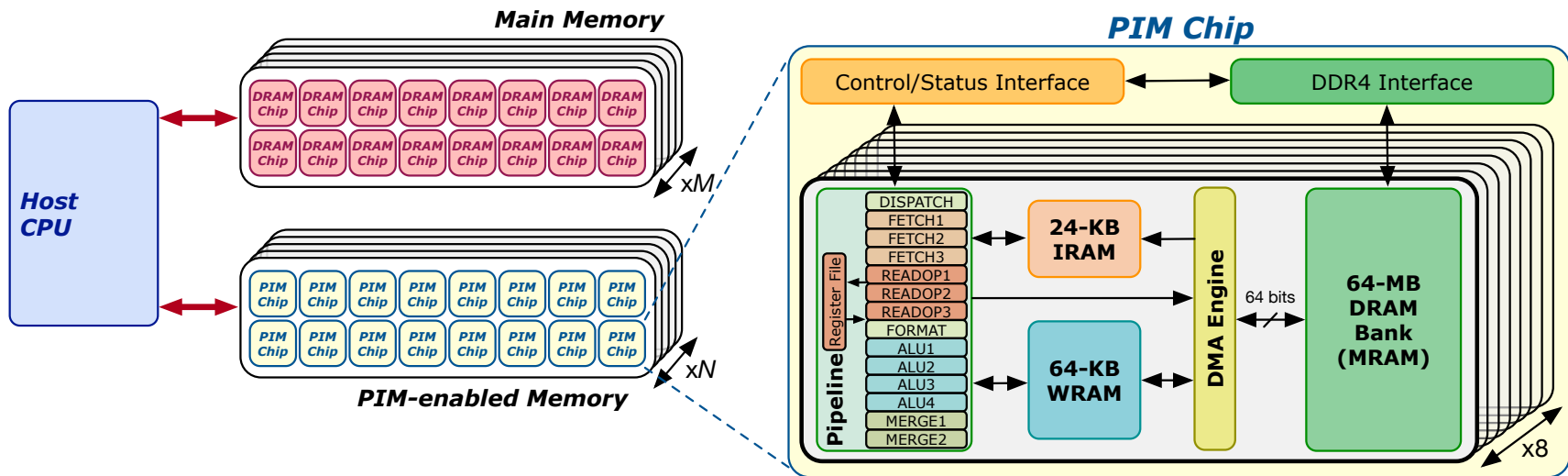
UPMEM Processing-in-DRAM Engine (2019)

- **Processing in DRAM Engine**
- Includes **standard DIMM modules**, with a **large number of DPU processors** combined with DRAM chips.
- Replaces **standard DIMMs**
 - DDR4 R-DIMM modules
 - 8GB+128 DPUs (16 PIM chips)
 - Standard 2x-nm DRAM process
 - **Large amounts of** compute & memory bandwidth



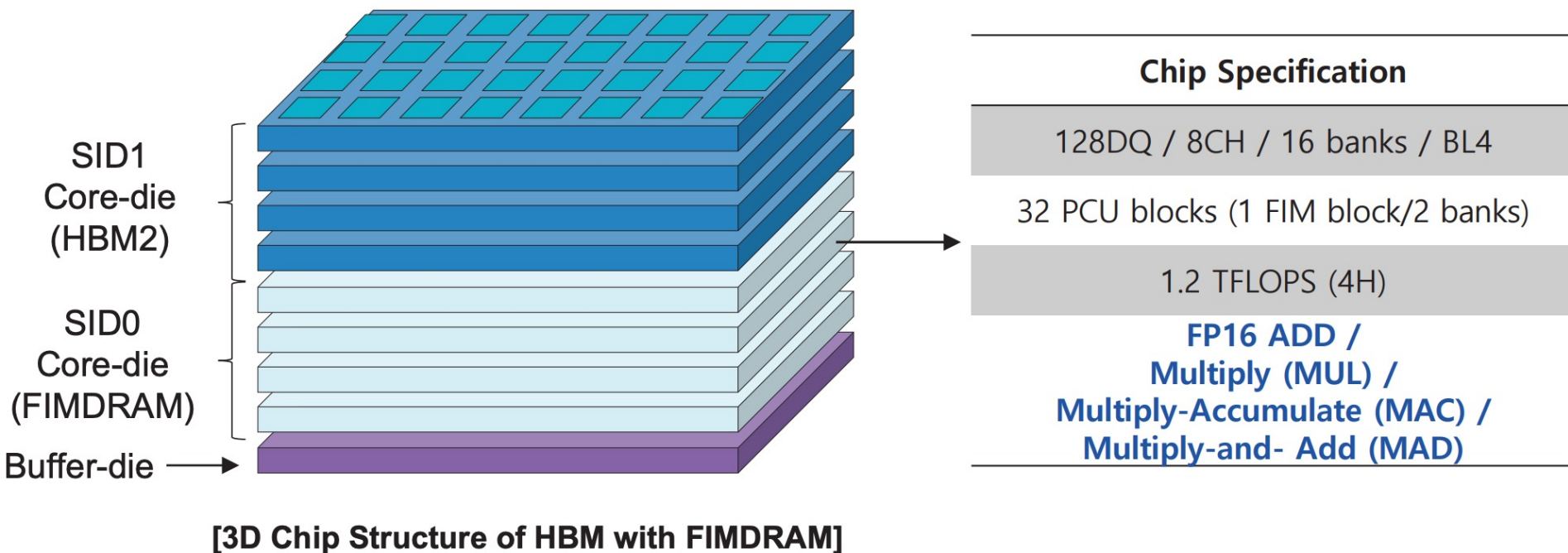
Recall: UPMEM PIM System Organization

- A UPMEM DIMM contains 8 or 16 chips
 - Thus, 1 or 2 ranks of 8 chips each
- Inside each PIM chip there are:
 - 8 64MB banks per chip: Main RAM (MRAM) banks
 - 8 DRAM Processing Units (DPUs) in each chip, 64 DPUs per rank



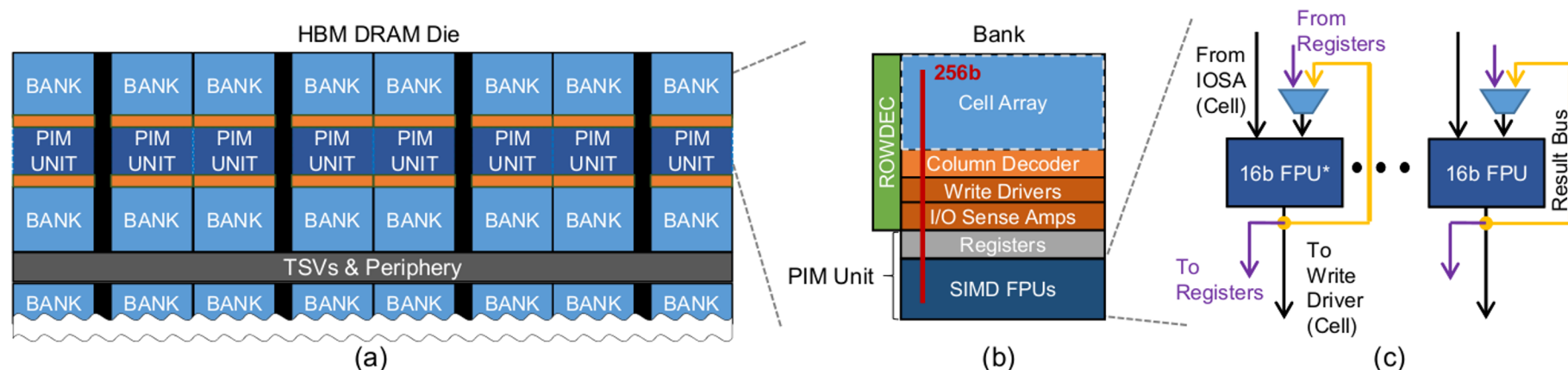
FIMDRAM: Chip Structure

■ FIMDRAM based on HBM2



FIMDRAM: System Organization (III)

- PIM units respond to standard DRAM column commands (RD or WR)
 - Compliant with unmodified JEDEC controllers
- They execute one wide-SIMD operation commanded by a PIM instruction with deterministic latency in a lock-step manner
- A PIM unit can get 16 16-bit operands from IOSAs, a register, and/or the result bus



SK Hynix Accelerator-in-Memory

SK Hynix Accelerator-in-Memory (2022)

SK hynix Develops PIM, Next-Generation AI Accelerator

February 16, 2022



Seoul, February 16, 2022

SK hynix (or “the Company”, www.skhynix.com) announced on February 16 that it has developed PIM*, a next-generation memory chip with computing capabilities.

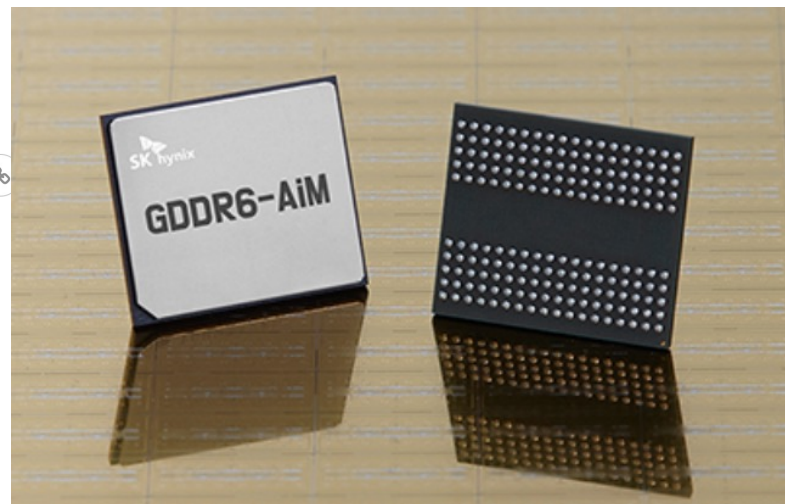
**PIM(Processing In Memory): A next-generation technology that provides a solution for data congestion issues for AI and big data by adding computational functions to semiconductor memory*

It has been generally accepted that memory chips store data and CPU or GPU, like human brain, process data. SK hynix, following its challenge to such notion and efforts to pursue innovation in the next-generation smart memory, has found a breakthrough solution with the development of the latest technology.

SK hynix plans to showcase its PIM development at the world’s most prestigious semiconductor conference, 2022 ISSCC*, in San Francisco at the end of this month. The company expects continued efforts for innovation of this technology to bring the memory-centric computing, in which semiconductor memory plays a central role, a step closer to the reality in devices such as smartphones.

**ISSCC: The International Solid-State Circuits Conference will be held virtually from Feb. 20 to Feb. 24 this year with a theme of “Intelligent Silicon for a Sustainable World”*

For the first product that adopts the PIM technology, SK hynix has developed a sample of GDDR6-AiM (Accelerator* in memory). The GDDR6-AiM adds computational functions to GDDR6* memory chips, which process data at 16Gbps. A combination of GDDR6-AiM with CPU or GPU instead of a typical DRAM makes certain computation speed 16 times faster. GDDR6-AiM is widely expected to be adopted for machine learning, high-performance computing, and big data computation and storage.



11.1 A 1nm 1.25V 8Gb, 16Gb/s/pin GDDR6-based Accelerator-in-Memory supporting 1TFLOPS MAC Operation and Various Activation Functions for Deep-Learning Applications

Seongju Lee, SK hynix, Icheon, Korea

In Paper 11.1, SK Hynix describes an 1nm, GDDR6-based accelerator-in-memory with a command set for deep-learning operation. The 8Gb design achieves a peak throughput of 1TFLOPS with 1GHz MAC operations and supports major activation functions to improve accuracy.

Accelerator-in-Memory (ISSCC 2022)

ISSCC 2022 / SESSION 11 / COMPUTE-IN-MEMORY AND

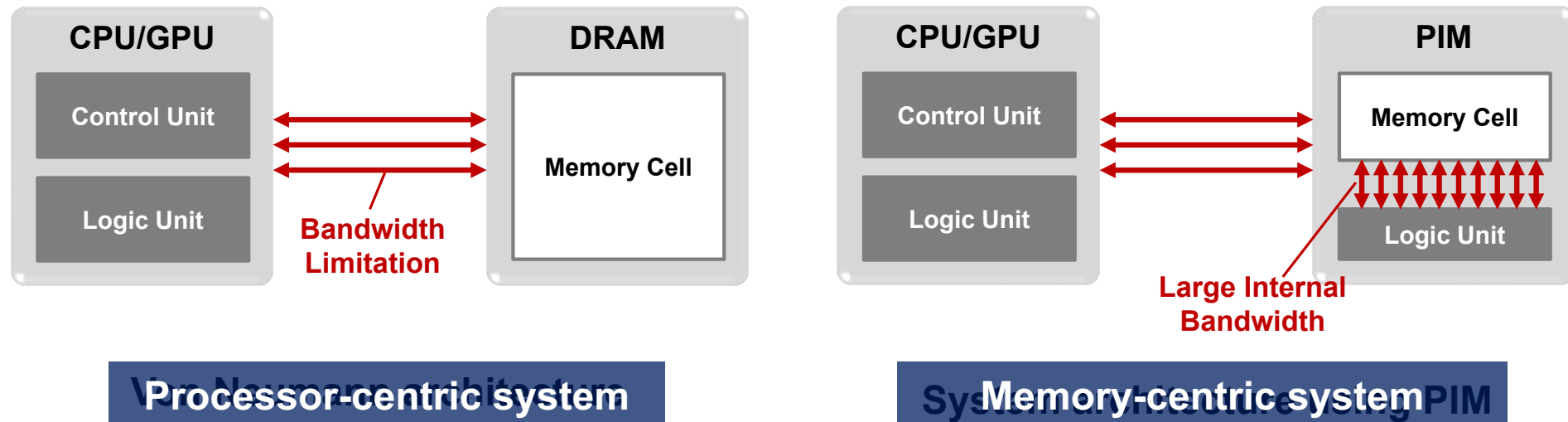
11.1 A 1ynm 1.25V 8Gb, 16Gb/s/pin GDDR6-based Accelerator-in-Memory supporting 1TFLOPS MAC Operation and Various Activation Functions for Deep-Learning Applications

Seongju Lee, Kyuyoung Kim, Sanghoon Oh, Joonhong Park, Gimoon Hong, Dongyoon Ka, Kyudong Hwang, Jeongje Park, Kyeongpil Kang, Jungyeon Kim, Junyeol Jeon, Nahsung Kim, Yongkee Kwon, Kornijcuk Vladimir, Woojae Shin, Jongsoon Won, Minkyu Lee, Hyunha Joo, Haerang Choi, Jaewook Lee, Donguc Ko, Younggun Jun, Keewon Cho, Ilwoong Kim, Choungki Song, Chunseok Jeong, Daehan Kwon, Jieun Jang, Il Park, Junhyun Chun, Joohwan Cho

SK hynix, Icheon, Korea

AiM: Exploiting Bank Parallelism

- Memory bandwidth is not enough for many ML workloads

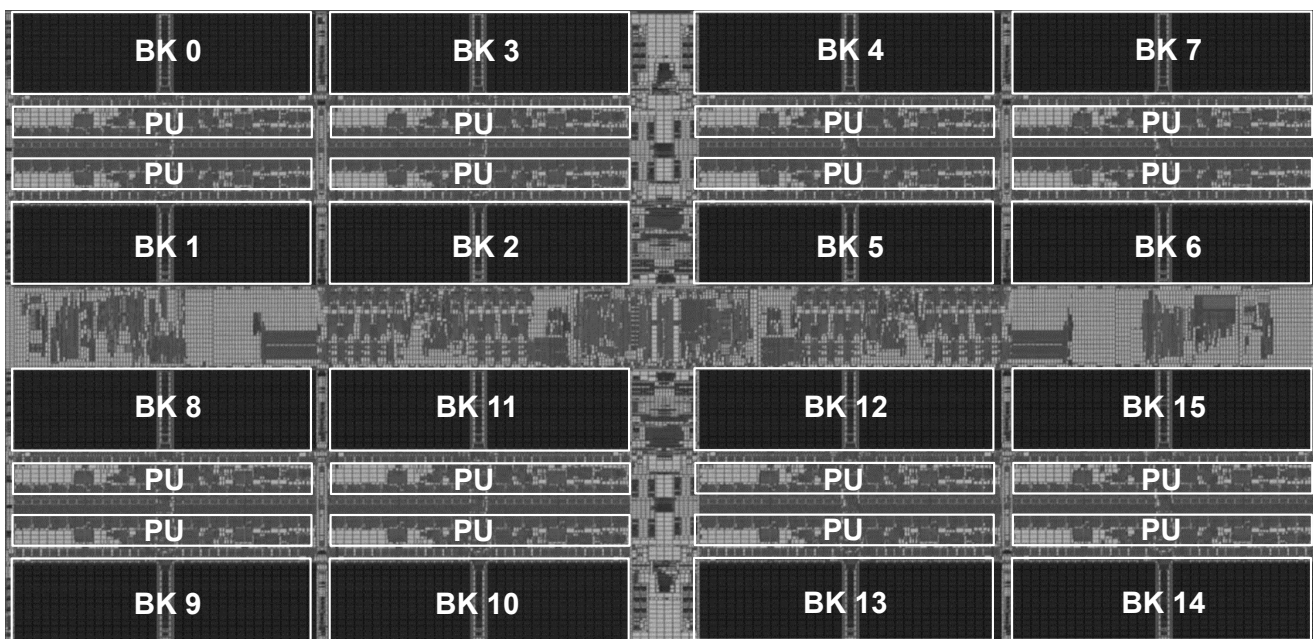


Onur Mutlu, ["Memory-Centric Computing"](#), Keynote Talk at the Thoughtworks Engineering for Research Symposium (E4R), Virtual, 19 February 2022.

AiM: Chip Implementation

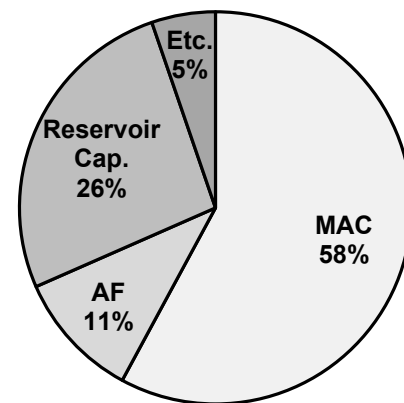
- 4 Gb AiM die with 16 processing units (PUs)

AiM Die Photograph



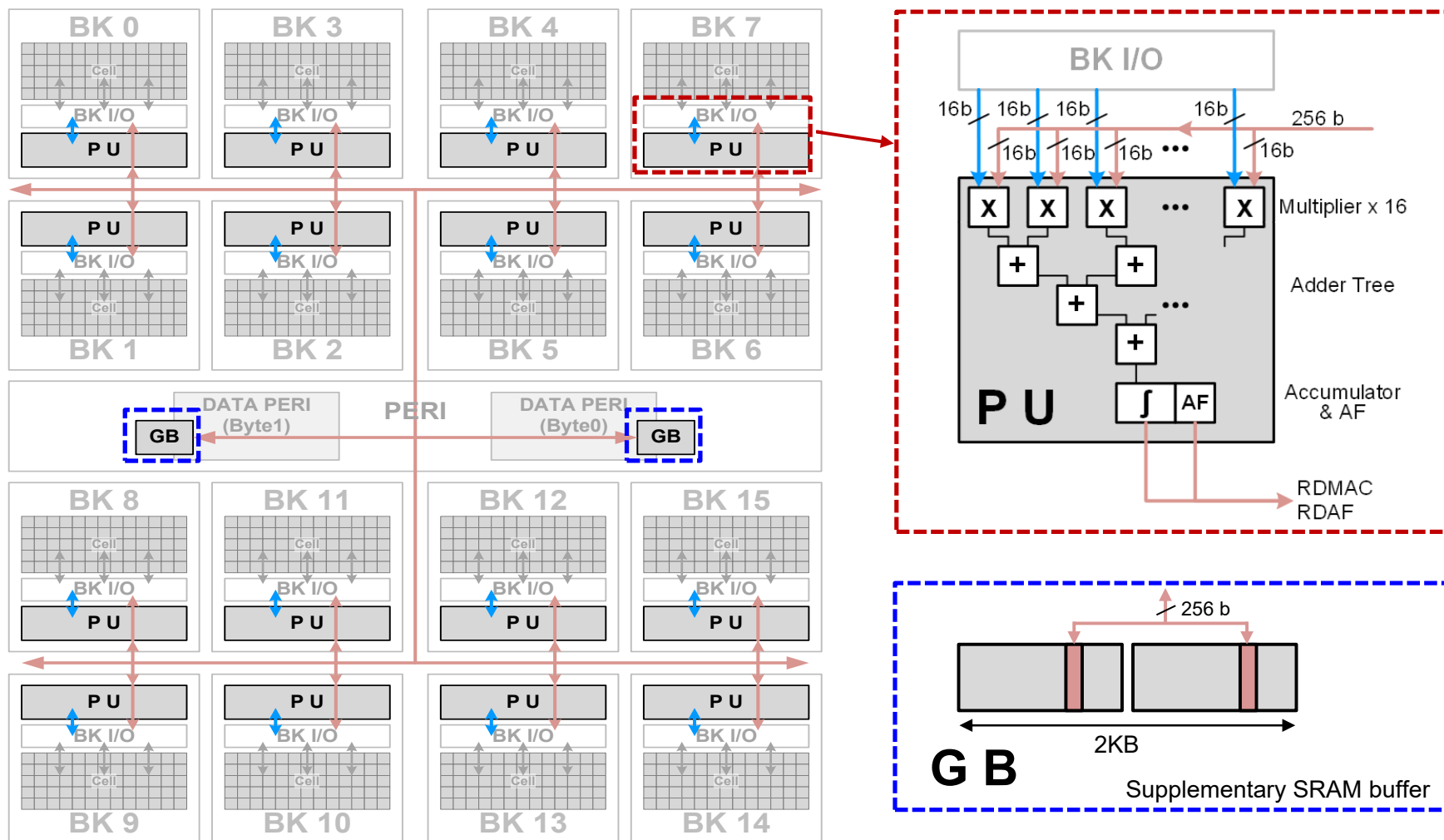
1 Process Unit (PU) Area

Total	0.19mm ²
MAC	0.11mm ²
Activation Function (AF)	0.02mm ²
Reservoir Cap.	0.05mm ²
Etc.	0.01mm ²



AiM: System Organization

■ GDDR6-based AiM architecture



AiM Commands

AiM: Command Set

■ New commands for computation

Type	CMD	Description
Bank Activation	ACT4, ACT16	Activate four/sixteen banks in parallel
	ACTAF4, ACTAF16	Activate rows storing activation function LUTs in four/sixteen banks in parallel
Compute	MACSB, MAC4B, MACAB	Perform MAC in one/four/sixteen banks in parallel
	AF	Compute activation function in all banks
	EWMUL	Perform element-wise multiplication
Data	WRBK	Write to all activated banks in parallel
	WRGB	Write to Global Buffer
	RDMAC	Read from MAC result register
	RDAF	Read from activation function result register
	WRBIAS	Write bias to MAC result register
	RDCP	Copy data from a bank to Global Buffer
	WRCP	Copy data from Global Buffer to a bank

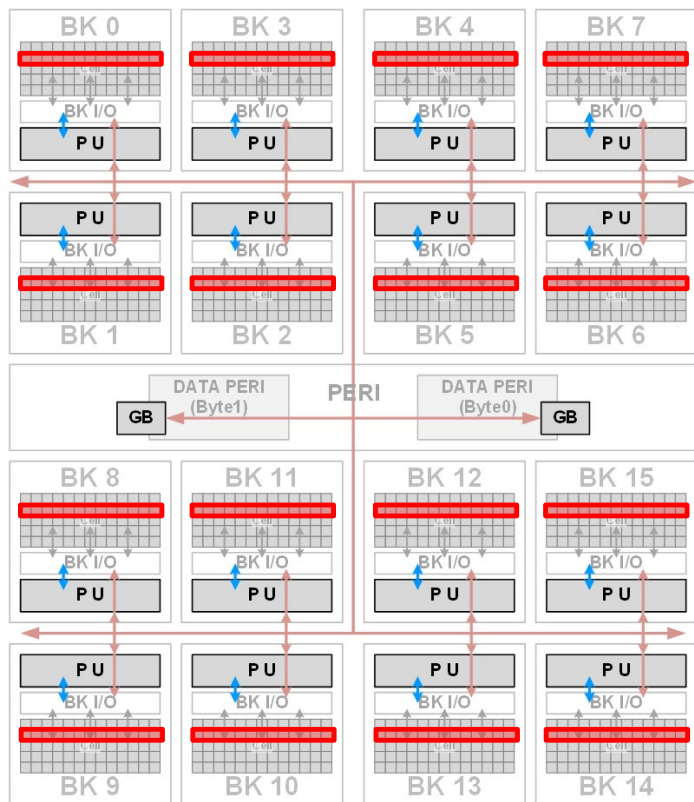
AiM: Command Set: ACT4, ACT16

- Activate 4 or 16 banks at once

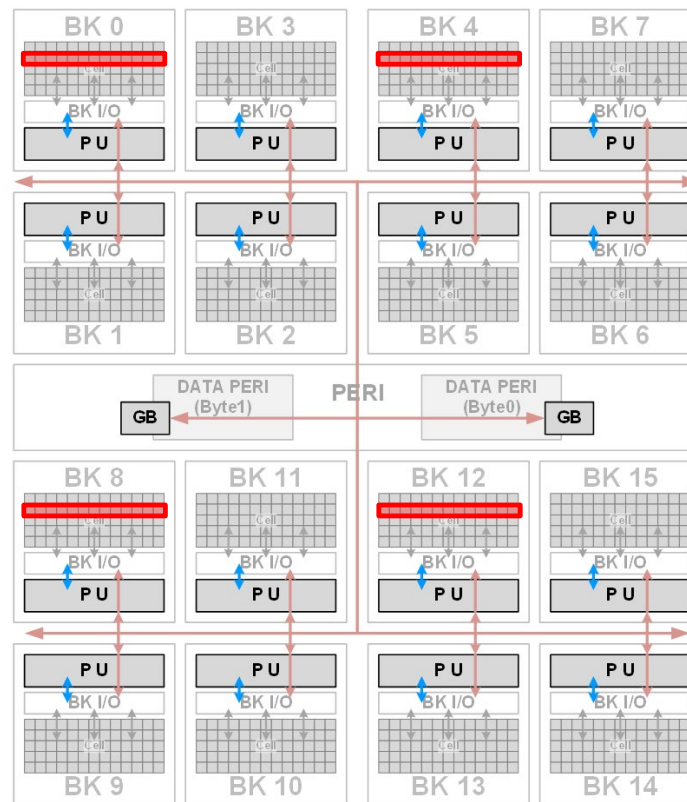
Activated Word Line

Type	CMD	Description
Bank Activation	ACT4, ACT16	Activate four/sixteen banks in parallel
	ACTAF4, ACTAF16	Activate rows storing activation function LUTs in four/sixteen banks in parallel
Compute	MACSB, MAC4B, MACAB	Perform MAC in one/four/sixteen banks in parallel
	AF	Compute activation function in all banks
	EWMUL	Perform element-wise multiplication
	WRBK	Write to all activated banks in parallel
Data	WRGB	Write to Global Buffer
	RDMAC	Read from MAC result register
	RDAF	Read from activation function result register
	WRBIAS	Write bias to MAC result register
	RDCP	Copy data from a bank to Global Buffer
	WRCP	Copy data from Global Buffer to a bank

16-bank Active



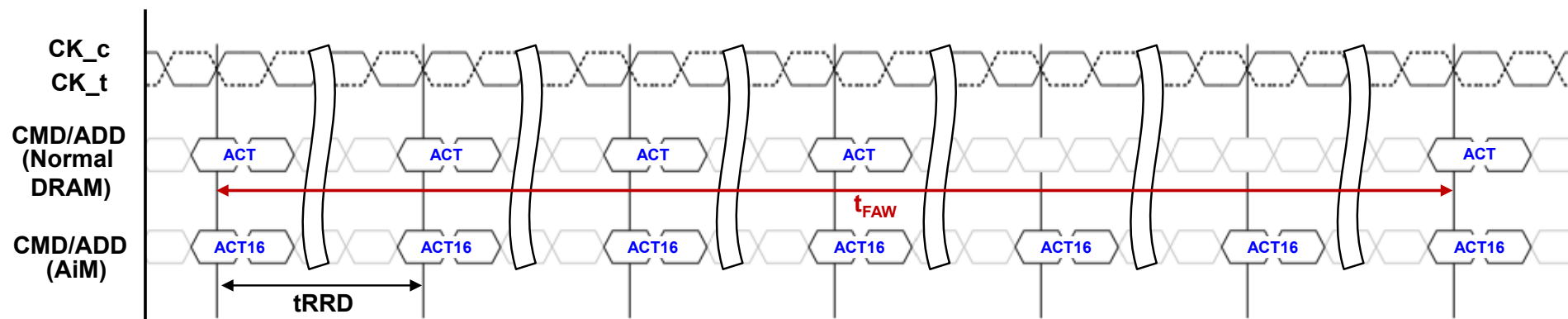
4-bank Active



AiM: Four Active Window (t_{FAW})

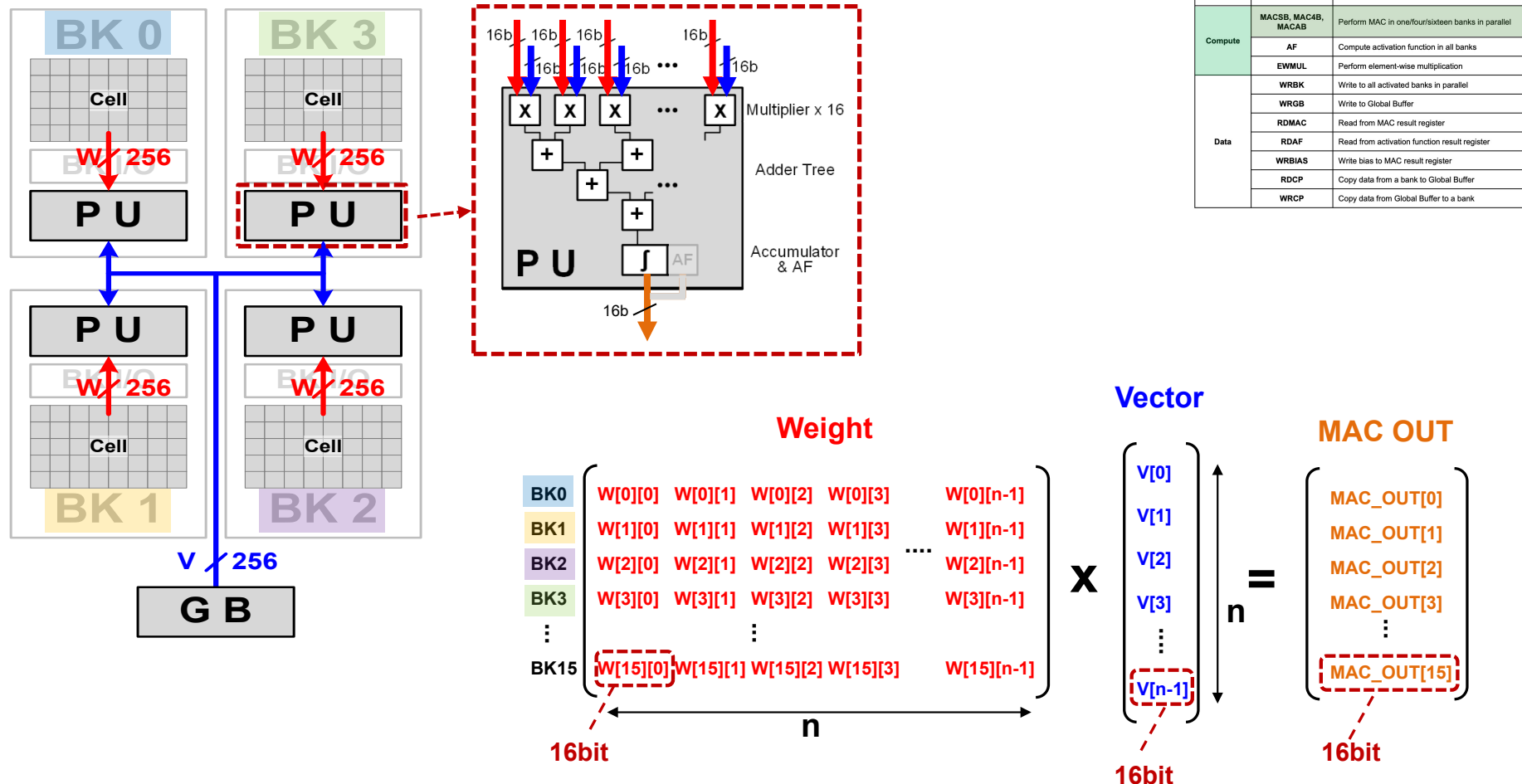
- ACT16 operations are possible without t_{FAW} constraint
 - In normal DRAM only 4 ACT are possible every t_{FAW}
 - Reservoir capacitor in each PU

Type	CMD	Description
Bank Activation	ACT4, ACT16	Activate four/sixteen banks in parallel
	ACTAF4, ACTAF16	Activate rows storing activation function LUTs in four/sixteen banks in parallel
Compute	MACSB, MAC4B, MACAB	Perform MAC in one/four/sixteen banks in parallel
	AF	Compute activation function in all banks
	EWMUL	Perform element-wise multiplication
Data	WRBK	Write to all activated banks in parallel
	WRGB	Write to Global Buffer
	RDMAC	Read from MAC result register
	RDAF	Read from activation function result register
	WRBIAS	Write bias to MAC result register
	RDCP	Copy data from a bank to Global Buffer
	WRCP	Copy data from Global Buffer to a bank



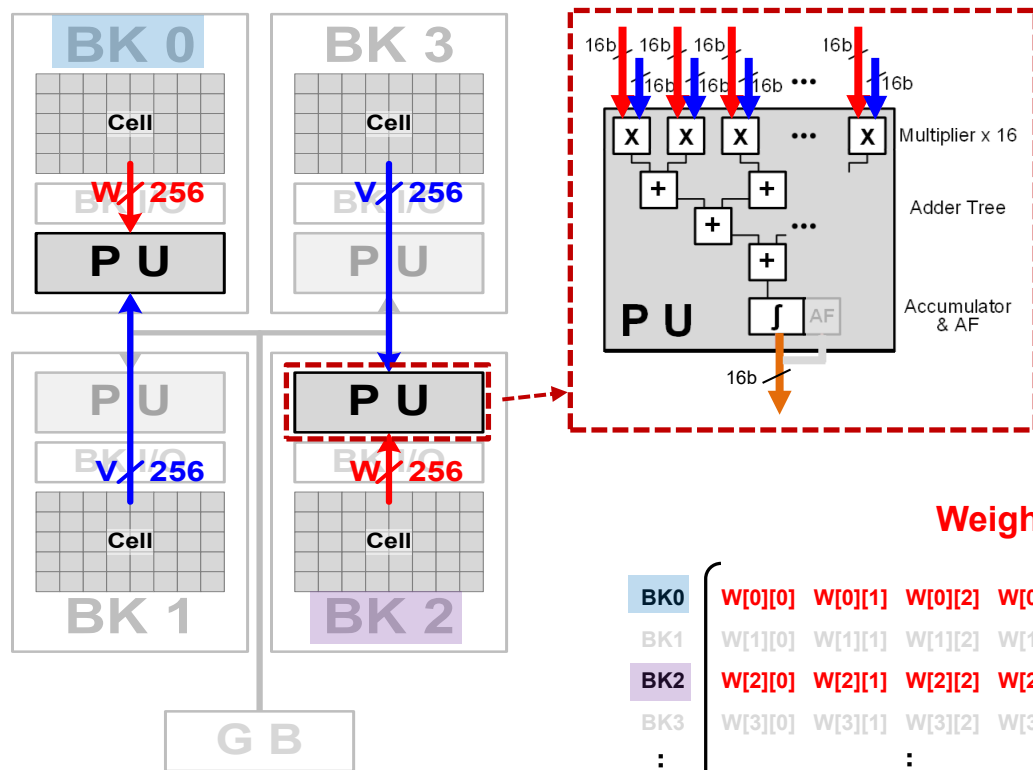
AiM: MAC Operation with Global Buffer

- MAC operation: **Weights from the banks**, **vectors from the GB**
- Operates with 16, 4, or 1 PU



AiM: MAC Operation without Global Buffer

- MAC operation: **Weights** and **vectors** from the banks
- Operates with 8, 4, or 1 PU



Type	CMD	Description
Bank Activation	ACT4, ACT16	Activate four/sixteen banks in parallel
	ACTAF4, ACTAF16	Activate rows storing activation function LUTs in four/sixteen banks in parallel
Compute	MACSB, MAC4B, MACAB	Perform MAC in one/four/sixteen banks in parallel
	AF	Compute activation function in all banks
Data	EWMUL	Perform element-wise multiplication
	WRBK	Write to all activated banks in parallel
	WRGB	Write to Global Buffer
	RDMAC	Read from MAC result register
	RDAF	Read from activation function result register
	WRBIAS	Write bias to MAC result register
	RDCP	Copy data from a bank to Global Buffer
	WRCP	Copy data from Global Buffer to a bank

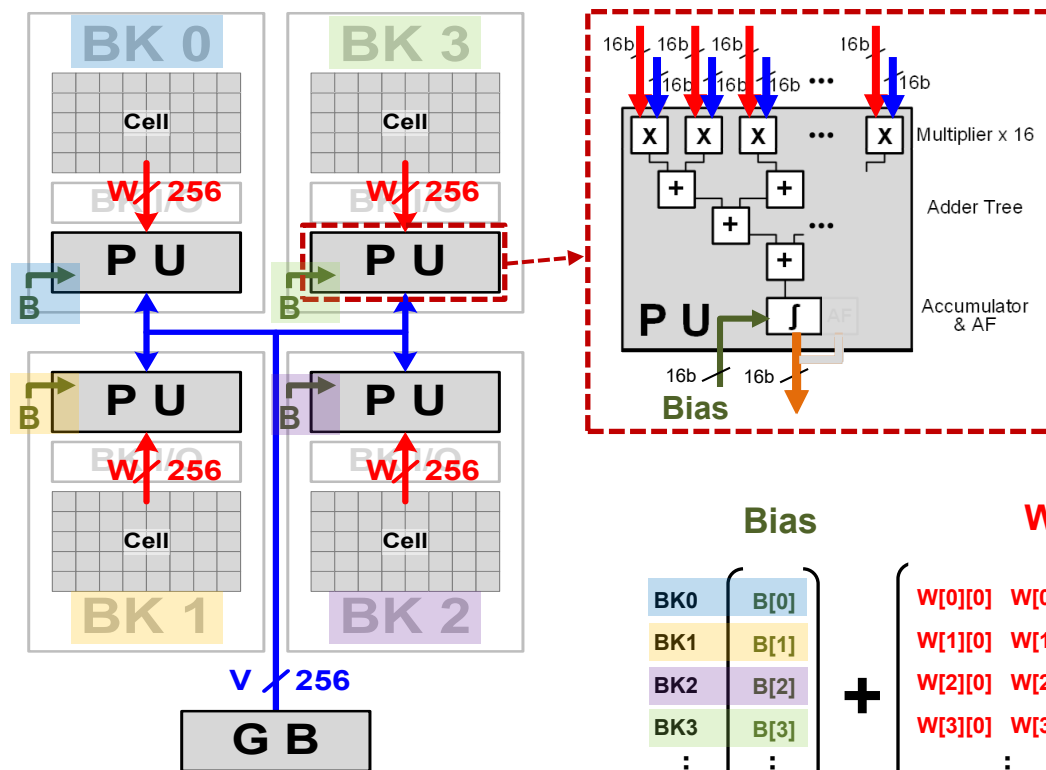
$$\begin{matrix}
 \text{Weight} & & \text{Vector} & & \text{MAC OUT} \\
 \begin{bmatrix}
 \text{BK0} & W[0][0] & W[0][1] & W[0][2] & W[0][3] & \dots & W[0][n-1] \\
 \text{BK1} & W[1][0] & W[1][1] & W[1][2] & W[1][3] & \dots & W[1][n-1] \\
 \text{BK2} & W[2][0] & W[2][1] & W[2][2] & W[2][3] & \dots & W[2][n-1] \\
 \text{BK3} & W[3][0] & W[3][1] & W[3][2] & W[3][3] & \dots & W[3][n-1] \\
 \vdots & \vdots & \vdots & \vdots & \vdots & \ddots & \vdots \\
 \text{BK15} & W[15][0] & W[15][1] & W[15][2] & W[15][3] & \dots & W[15][n-1]
 \end{bmatrix}
 & \mathbf{x} &
 \begin{bmatrix}
 V[0] \\
 V[1] \\
 V[2] \\
 V[3] \\
 \vdots \\
 V[n-1]
 \end{bmatrix}
 & = &
 \begin{bmatrix}
 \text{MAC_OUT}[0] \\
 \text{MAC_OUT}[1] \\
 \text{MAC_OUT}[2] \\
 \text{MAC_OUT}[3] \\
 \vdots \\
 \text{MAC_OUT}[15]
 \end{bmatrix}
 \end{matrix}$$

n (rows) n (columns) n (rows)

16bit (row) 16bit (column) 16bit (row)

AiM: Write Bias

- Biases can be added to MAC results
 - Different biases in 16 banks at the same time



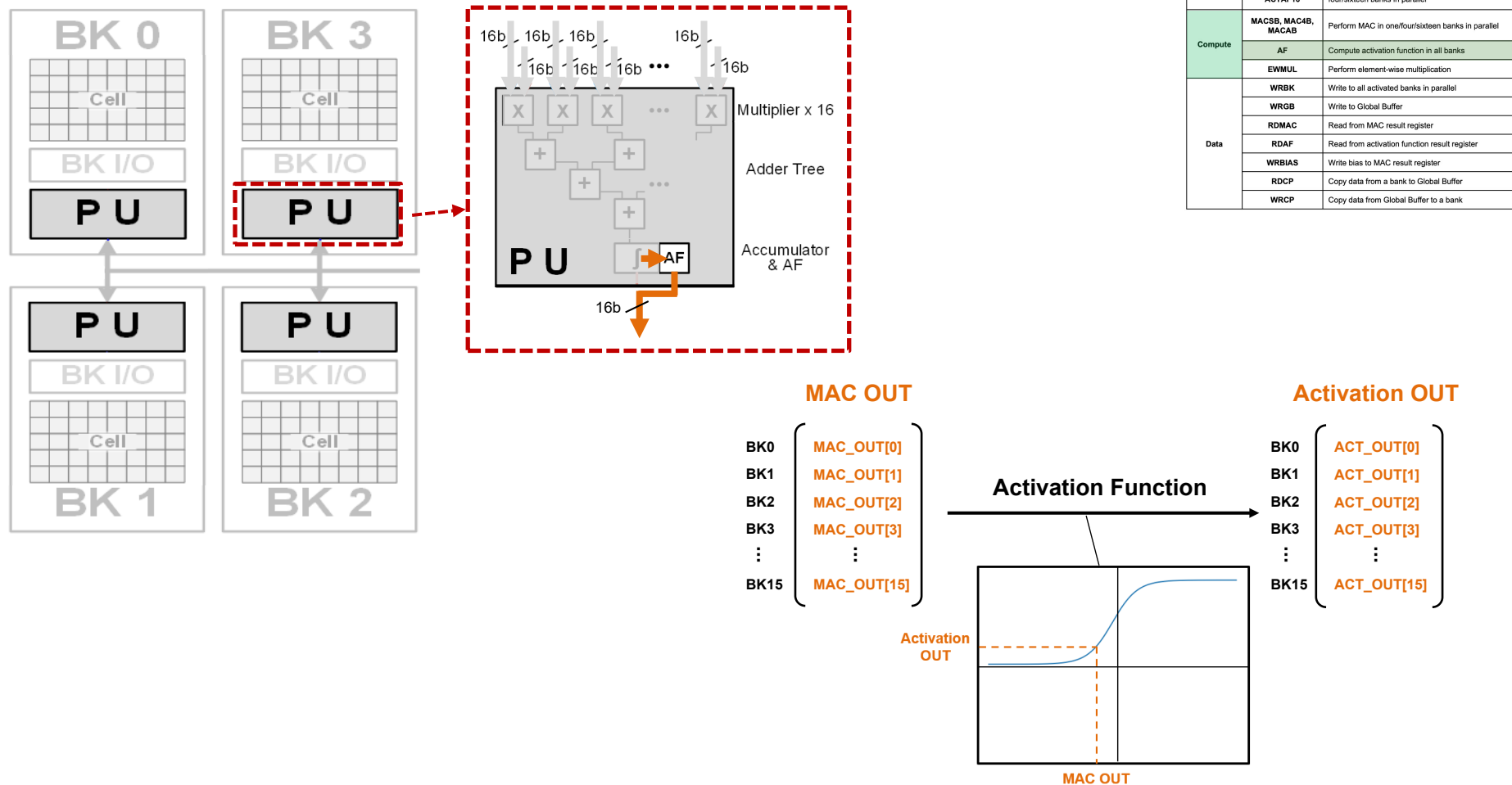
Type	CMD	Description
Bank Activation	ACT4, ACT16	Activate four/sixteen banks in parallel
	ACTAF4, ACTAF16	Activate rows storing activation function LUTs in four/sixteen banks in parallel
Compute	MACSB, MAC4B, MACAB	Perform MAC in one/four/sixteen banks in parallel
	AF	Compute activation function in all banks
	EWMUL	Perform element-wise multiplication
Data	WRBK	Write to all activated banks in parallel
	WRGB	Write to Global Buffer
	RDMAC	Read from MAC result register
	RDAF	Read from activation function result register
	WRBIAS	Write bias to MAC result register
	RDCP	Copy data from a bank to Global Buffer
	WRCP	Copy data from Global Buffer to a bank

$$\begin{array}{c}
 \text{Bias} \\
 \begin{array}{|c|} \hline \text{BK0} \\ \hline \text{BK1} \\ \hline \text{BK2} \\ \hline \text{BK3} \\ \hline \vdots \\ \hline \text{BK15} \\ \hline \end{array}
 \begin{array}{|c|} \hline \text{B}[0] \\ \hline \text{B}[1] \\ \hline \text{B}[2] \\ \hline \text{B}[3] \\ \hline \vdots \\ \hline \text{B}[15] \\ \hline \end{array}
 \end{array}
 +
 \begin{array}{|c|} \hline \text{Weight} \\ \hline \begin{array}{ccc} \text{W}[0][0] & \text{W}[0][1] & \text{W}[0][n-1] \\ \text{W}[1][0] & \text{W}[1][1] & \text{W}[1][n-1] \\ \vdots & \vdots & \vdots \\ \text{W}[15][0] & \text{W}[15][1] & \text{W}[15][n-1] \end{array} \\ \hline \end{array}
 \times
 \begin{array}{|c|} \hline \text{Vector} \\ \hline \begin{array}{c} \text{V}[0] \\ \text{V}[1] \\ \text{V}[2] \\ \text{V}[3] \\ \vdots \\ \text{V}[n-1] \end{array} \\ \hline \end{array}
 =
 \begin{array}{|c|} \hline \text{MAC OUT} \\ \hline \begin{array}{c} \text{MAC_OUT}[0] \\ \text{MAC_OUT}[1] \\ \text{MAC_OUT}[2] \\ \text{MAC_OUT}[3] \\ \vdots \\ \text{MAC_OUT}[15] \end{array} \\ \hline \end{array}$$

16bit
16bit
n
16bit
16bit

AiM: Activation Function

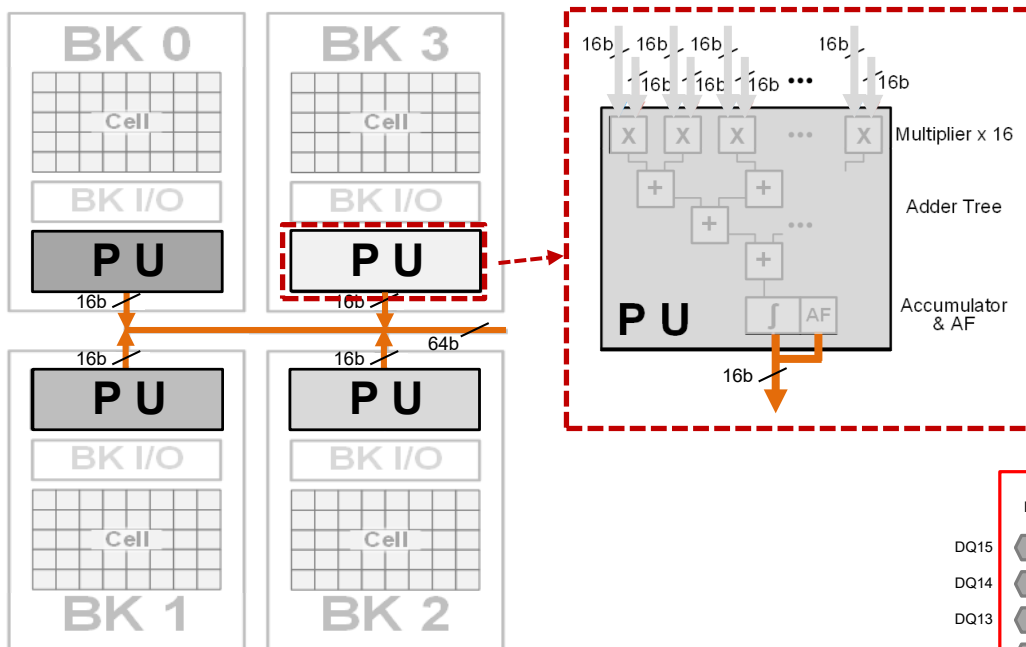
■ Activation function (AF): MAC result as input



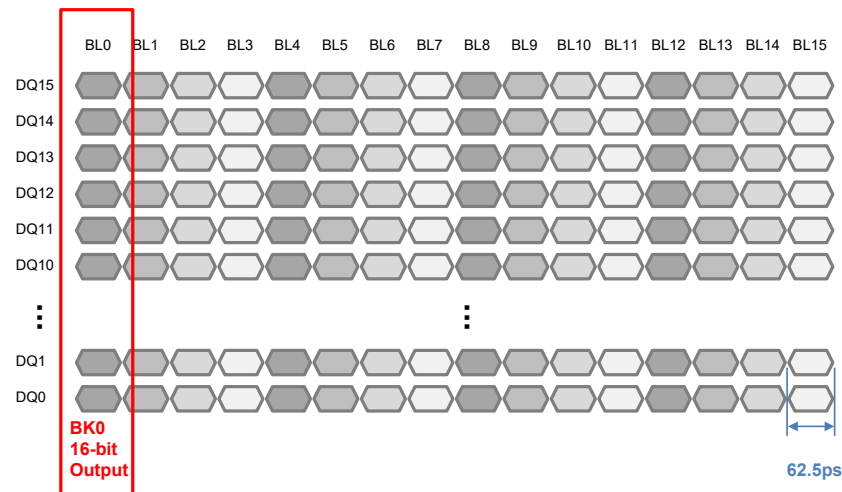
AiM: Read MAC / Read AF

■ Read outside DRAM

□ 16 banks x 16 bits = 256 bits

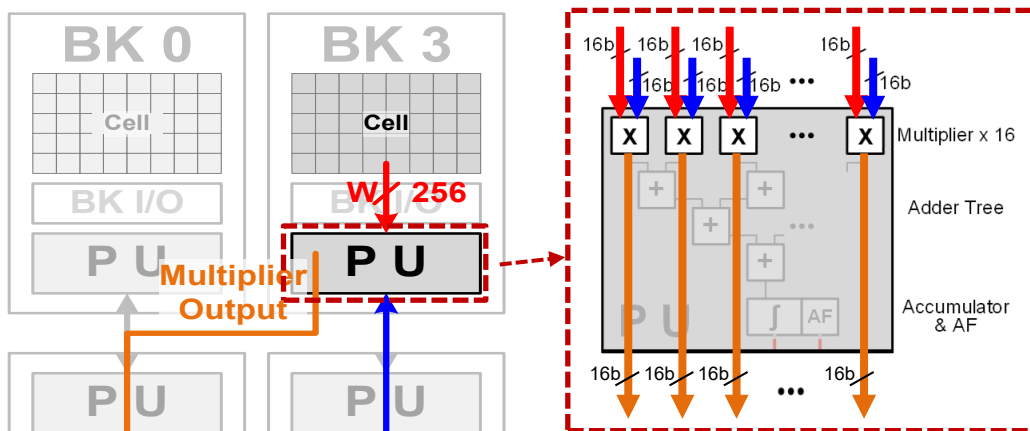


Type	CMD	Description
Bank Activation	ACT4, ACT16	Activate four/sixteen banks in parallel
	ACTAF4, ACTAF16	Activate rows storing activation function LUTs in four/sixteen banks in parallel
Compute	MACSB, MAC4B, MACAB	Perform MAC in one/four/sixteen banks in parallel
	AF	Compute activation function in all banks
	EWMUL	Perform element-wise multiplication
Data	WRBK	Write to all activated banks in parallel
	WRGB	Write to Global Buffer
	RDMAC	Read from MAC result register
	RDAF	Read from activation function result register
	WRBIAS	Write bias to MAC result register
	RDCP	Copy data from a bank to Global Buffer
	WRCP	Copy data from Global Buffer to a bank



AiM: Element-wise Multiplication

- Multiplies **weights** and **vectors**
 - The adder tree is disabled



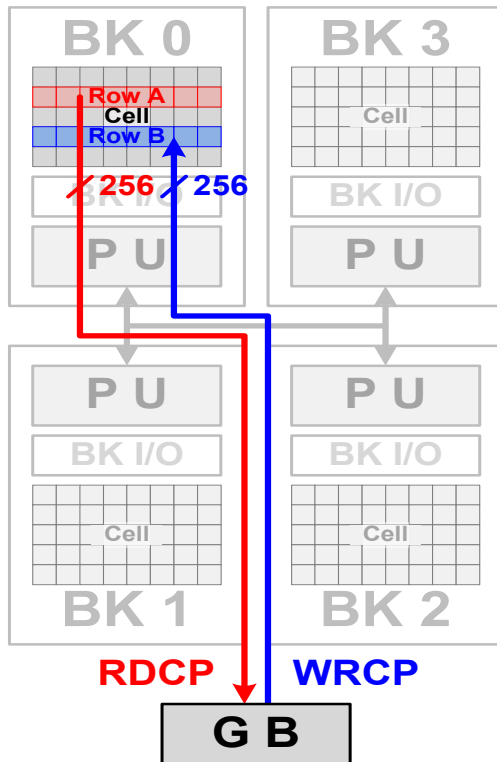
Type	CMD	Description
Bank Activation	ACT4, ACT16	Activate four/sixteen banks in parallel
	ACTAF4, ACTAF16	Activate rows storing activation function LUTs in four/sixteen banks in parallel
Compute	MACS8, MAC48, MACAB	Perform MAC in one/four/sixteen banks in parallel
	AF	Compute activation function in all banks
	EWMUL	Perform element-wise multiplication
Data	WRBK	Write to all activated banks in parallel
	WRGB	Write to Global Buffer
	RDMAC	Read from MAC result register
	RDAF	Read from activation function result register
	WRBIAS	Write bias to MAC result register
	RDCP	Copy data from a bank to Global Buffer
	WRCP	Copy data from Global Buffer to a bank

$$\begin{array}{c}
 \text{Weight} \\
 \begin{bmatrix} w[0] \\ w[1] \\ w[2] \\ w[3] \\ \vdots \\ w[n-1] \end{bmatrix} \\
 \text{16bit}
 \end{array}
 \times
 \begin{array}{c}
 \text{Vector} \\
 \begin{bmatrix} v[0] \\ v[1] \\ v[2] \\ v[3] \\ \vdots \\ v[n-1] \end{bmatrix} \\
 \text{16bit}
 \end{array}
 =
 \begin{array}{c}
 \text{Multiplier Output} \\
 \begin{bmatrix} \text{MUL_OUT}[0] \\ \text{MUL_OUT}[1] \\ \text{MUL_OUT}[2] \\ \text{MUL_OUT}[3] \\ \vdots \\ \text{MUL_OUT}[n-1] \end{bmatrix} \\
 \text{16bit}
 \end{array}$$

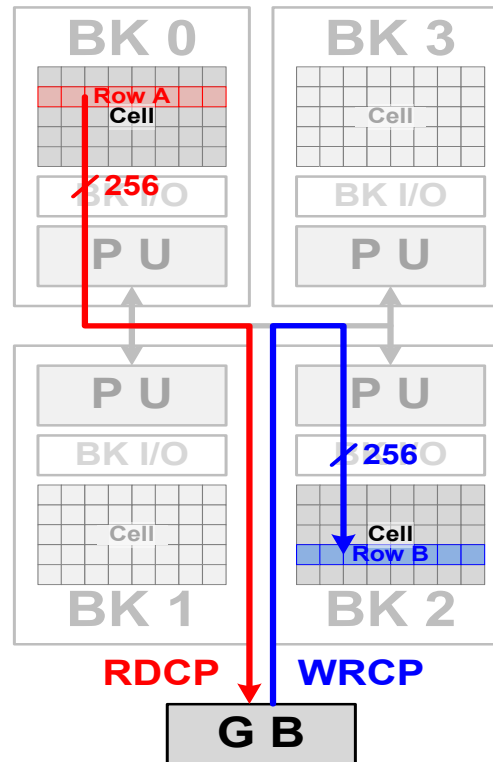
AiM: Copy Operation

- Copy operation using GB as a intermediate buffer
 - It copies 2 KB from one row to another row

Copy to Same Bank



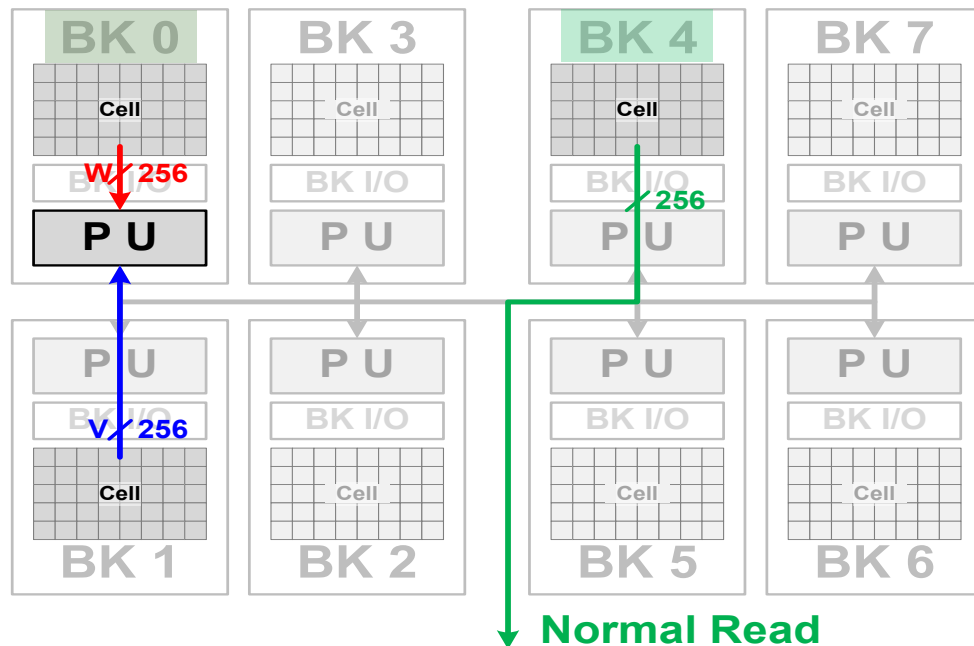
Copy to Another Bank



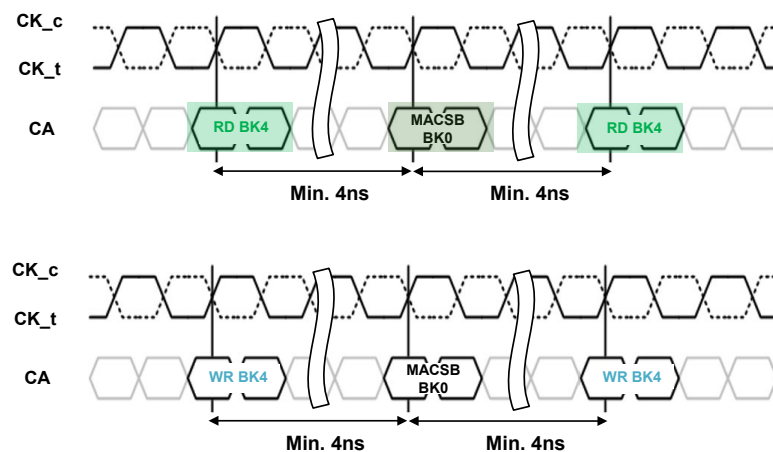
Type	CMD	Description
Bank Activation	ACT4, ACT16	Activate four/sixteen banks in parallel
	ACTAF4, ACTAF16	Activate rows storing activation function LUTs in four/sixteen banks in parallel
Compute	MACSB, MAC4B, MACAB	Perform MAC in one/four/sixteen banks in parallel
	AF	Compute activation function in all banks
	EWMUL	Perform element-wise multiplication
	WRBK	Write to all activated banks in parallel
Data	WRGB	Write to Global Buffer
	RDMAC	Read from MAC result register
	RDAF	Read from activation function result register
	WRBIAS	Write bias to MAC result register
	RDCP	Copy data from a bank to Global Buffer
	WRCP	Copy data from Global Buffer to a bank

AiM: Simultaneous Computation and Access

- Computation can be performed during normal read/write



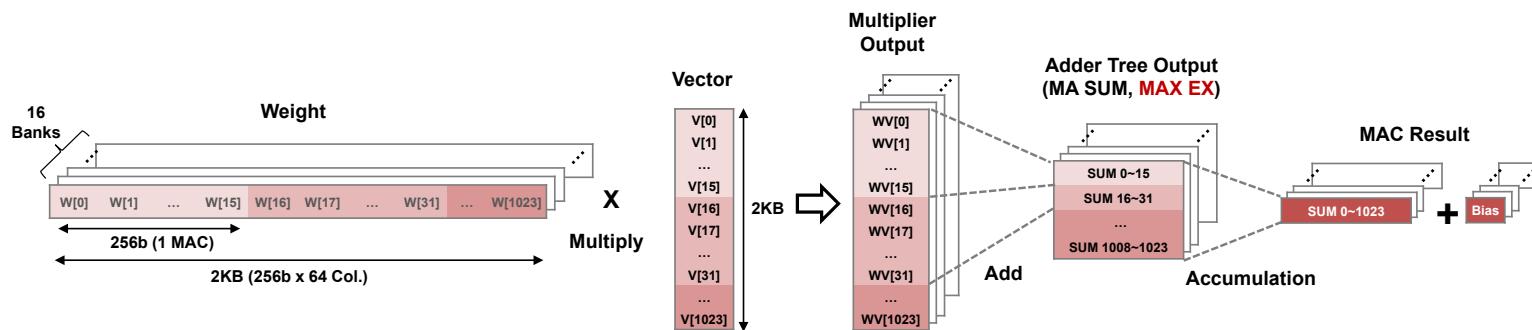
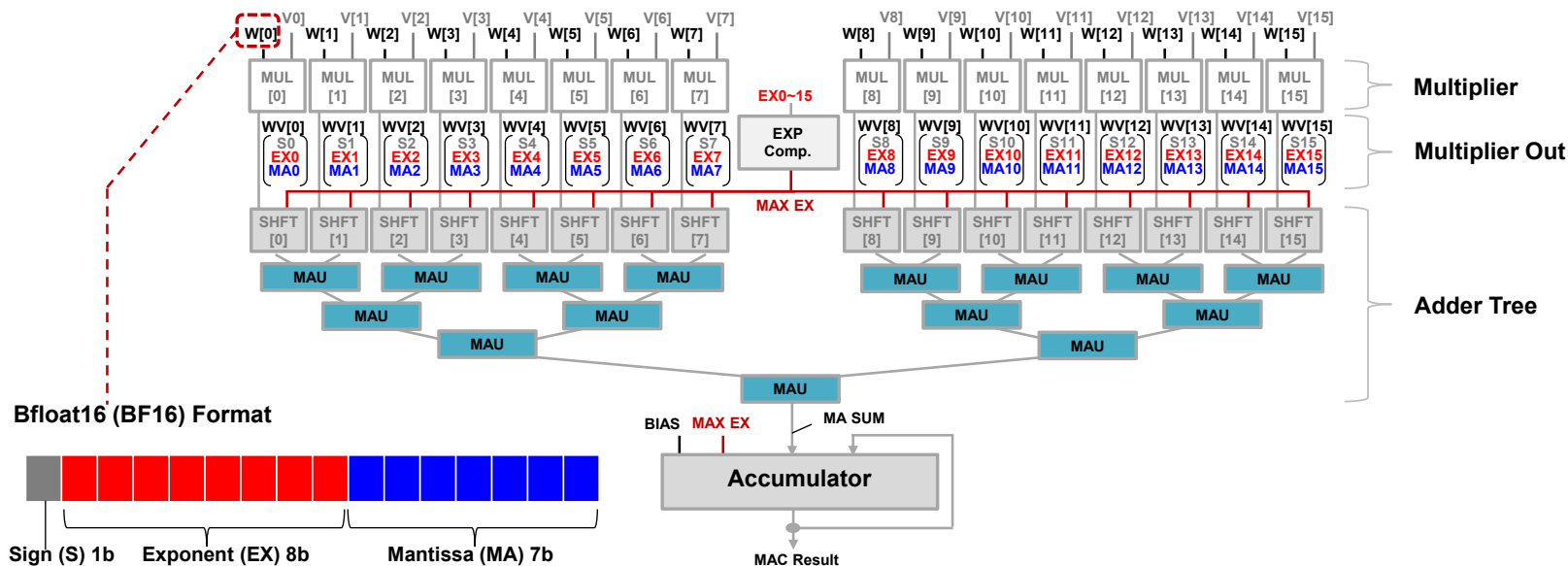
Type	CMD	Description
Bank Activation	ACT4, ACT16	Activate four/sixteen banks in parallel
	ACTAF4, ACTAF16	Activate rows storing activation function LUTs in four/sixteen banks in parallel
Compute	MACSB, MAC4B, MACAB	Perform MAC in one/four/sixteen banks in parallel
	AF	Compute activation function in all banks
	EWMUL	Perform element-wise multiplication
Data	WRBK	Write to all activated banks in parallel
	WRGB	Write to Global Buffer
	RDMAC	Read from MAC result register
	RDAF	Read from activation function result register
	WRBIAS	Write bias to MAC result register
	RDCP	Copy data from a bank to Global Buffer
	WRCP	Copy data from Global Buffer to a bank



AiM Microarchitecture

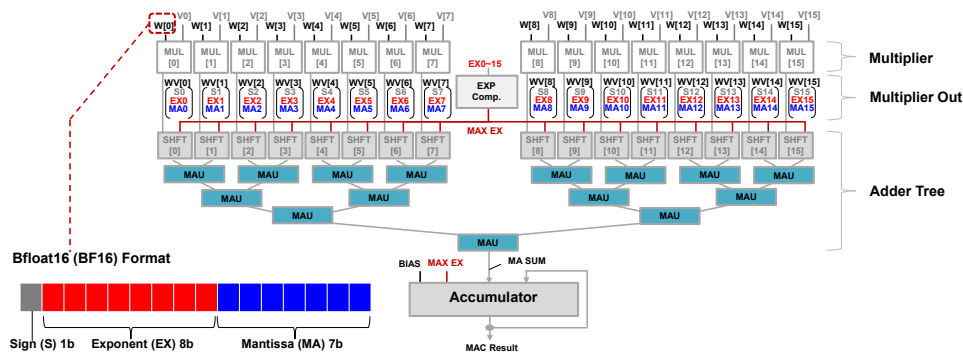
AiM: MAC Circuit

- 16 multipliers, adder tree, and accumulator
 - Bfloat16 (BF16) format

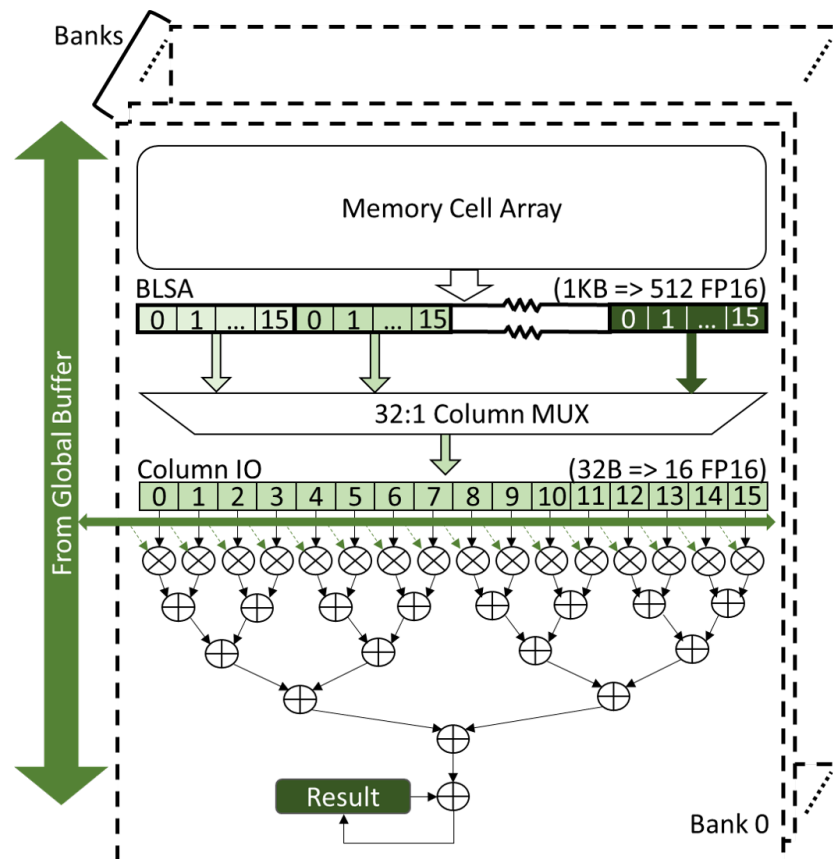


AiM: MAC Circuit: Prior Work (I)

- 16 multipliers, adder tree, and accumulator

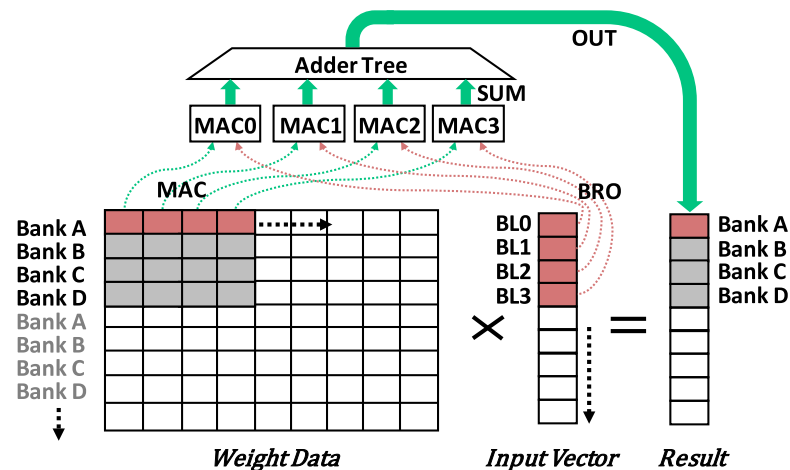
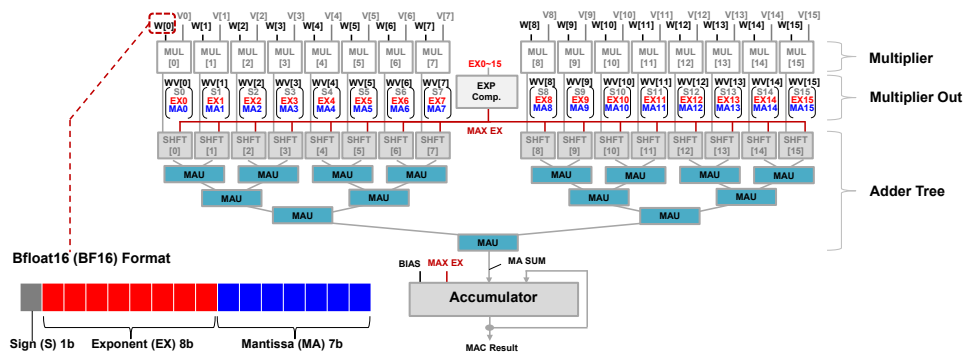


He et al, "Newton: A DRAM-maker's accelerator-in-memory (AiM) architecture for machine learning," MICRO 2020

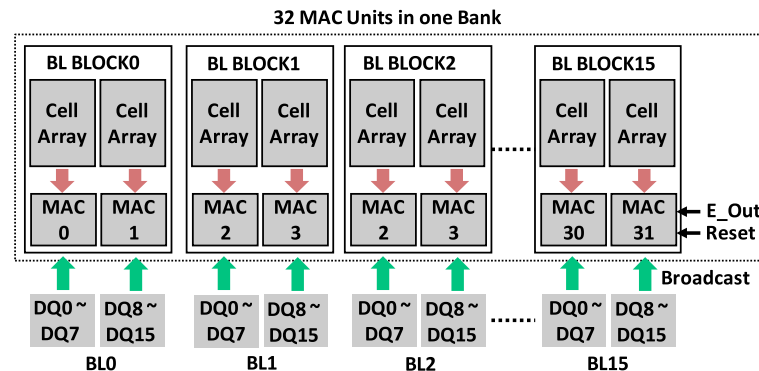


AiM: MAC Circuit: Prior Work (II)

- 16 multipliers, adder tree, and accumulator

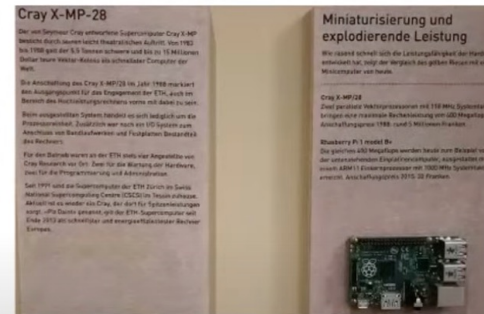


Shin et al, "McDRAM: Low latency and energy-efficient matrix computations in DRAM," IEEE TCADICS (2018)



Lecture on SIMD Processing

CRAY X-MP-28 @ ETH (CAB, E Floor)



DEPARTMENT OF INFORMATION TECHNOLOGY AND ELECTRICAL ENGINEERING (D-ITET)

Digital Design & Comp. Arch. - Lecture 20: SIMD Processing (Vector and Array Processors) (Spring'21)

2,677 views • Streamed live on May 14, 2021

69 1 SHARE SAVE ...



Onur Mutlu Lectures
19.2K subscribers

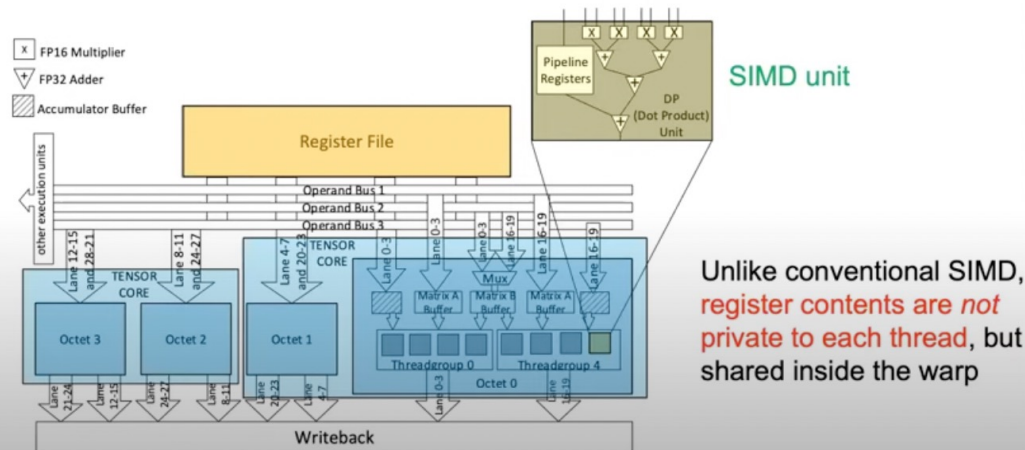
SUBSCRIBED



Lecture on SIMD Processing and GPUs

Tensor Core Microarchitecture (Volta)

- Each warp utilizes **two tensor cores**
- Each tensor core contains **two "octets"**
 - **16 SIMD units per tensor core** (8 per octet)
 - 4x4 matrix-multiply and accumulate each cycle per tensor core



~~Proposed* tensor core microarchitecture~~

HetSys Course: Lecture 2: SIMD Processing and GPUs (Spring 2022)

380 views • Premiered Mar 22, 2022

9 DISLIKE SHARE CLIP SAVE ...



Onur Mutlu Lectures

23.6K subscribers

Project & Seminar, ETH Zürich, Spring 2022

Hands-on Acceleration on Heterogeneous Computing Systems (

https://safari.ethz.ch/projects_and_s...)

SUBSCRIBED



Parallel Reduction on GPU

Divergence-Free Mapping (II)

- Program with high SIMD utilization

```
__shared__ float partialSum[]  
  
unsigned int t = threadIdx.x;  
  
for(int stride = blockDim.x; stride > 0; stride >> 1){  
  
    __syncthreads();  
  
    if (t < stride)  
        partialSum[t] += partialSum[t + stride];  
  
}
```

ETH zürich SAFARI

Juan Gomez L...

33:56 / 1:15:55

CC HD 2

Heterogeneous Systems Course: Meeting 6: Parallel Patterns: Reduction (Fall 2021)

411 views • Streamed live on Nov 11, 2021

👍 23 🗨️ 0 ➦ SHARE ➦ SAVE ...



Onur Mutlu Lectures

20.1K subscribers

SUBSCRIBED



Project & Seminar, ETH Zürich, Fall 2021

Hands-on Acceleration on Heterogeneous Computing Systems (

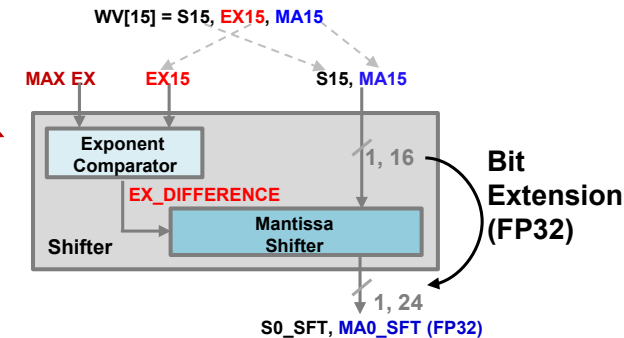
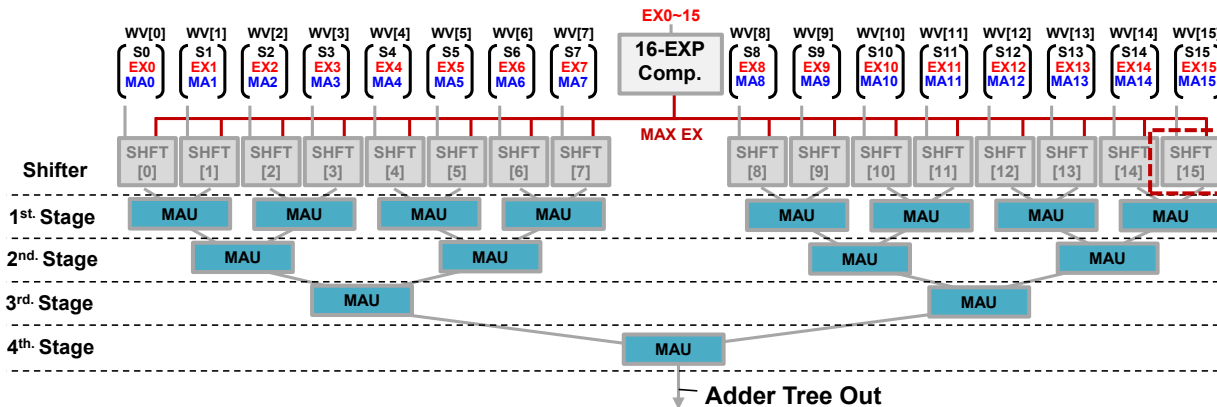
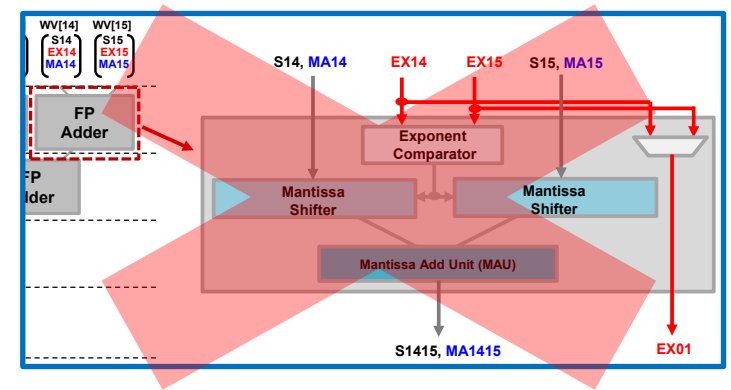
https://safari.ethz.ch/projects_and_s...)

<https://youtu.be/PN0h5XPhJIM>

AiM: Adder Tree: Bank-wide Mantisa Shift (I)

Bank-wide Mantisa Shift (BWMS)

- Find MAX EX of 16 EXs
- Obtain the differences
- Shift all MAs by the differences
- Perform MA additions

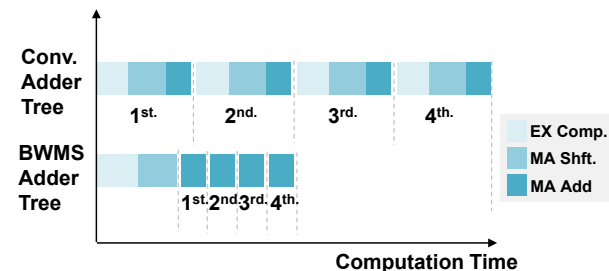
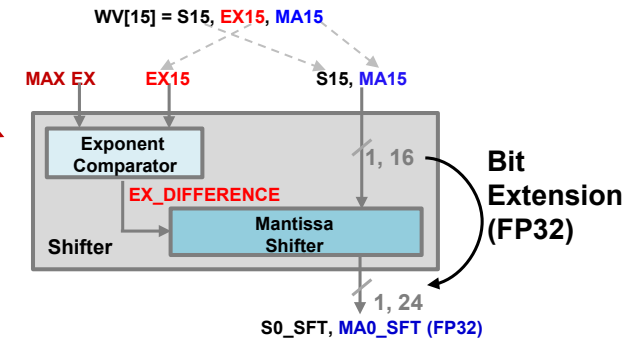
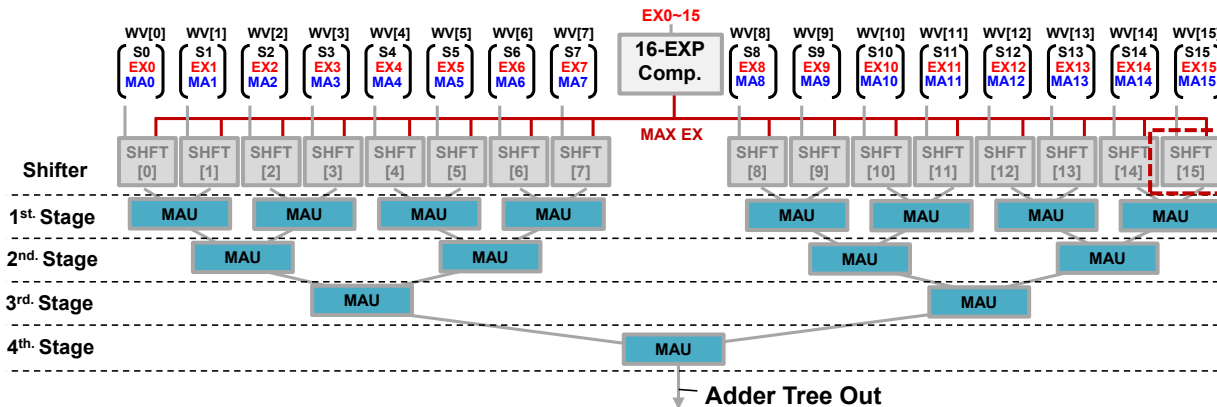
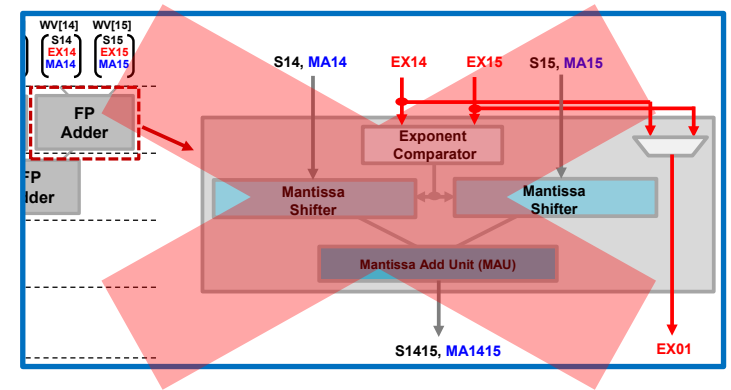


	Before Shifter			After Shifter	
MUL OUT	Exponent (EX) (8 bits)	Sign (S), Mantissa (MA) (1, 16 bits)	Diff. w/ Max EX	Exponent (EX) (8 bits)	Sign(S), Mantissa(MA) (1, 24 bits)
WV[0]	00001100	+1.111110011111011	3	00001111	+0.00111111001111101100000
WV[1]	00001111	+1.000000011111011	Max EX		+1.00000001111101100000000
WV[2]	00000100	-1.111110000000000	11		-1.111111111111110000000
⋮	⋮	⋮	⋮		⋮
WV[15]	00001000	+1.000000011111111	7		+0.00000010000000011111110

AiM: Adder Tree: Bank-wide Mantissa Shift (II)

Bank-wide Mantissa Shift (BWMS)

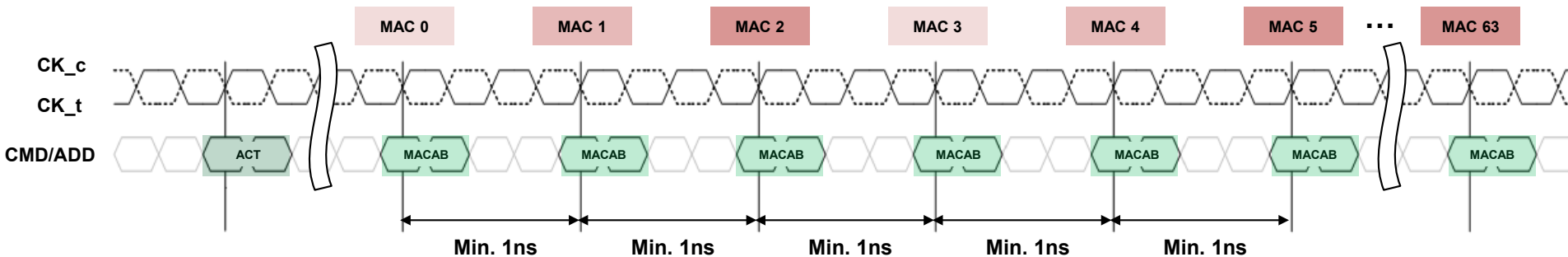
- Find MAX EX of 16 EXs
- Obtain the differences
- Shift all MAs by the differences
- Perform MA additions



AiM: Accumulation Operation

- Up to 64 MAC commands between ACT
 - 1 MAC every 1ns (t_{CCDS})
 - 64 x 256 bits = 2 KB (row size)

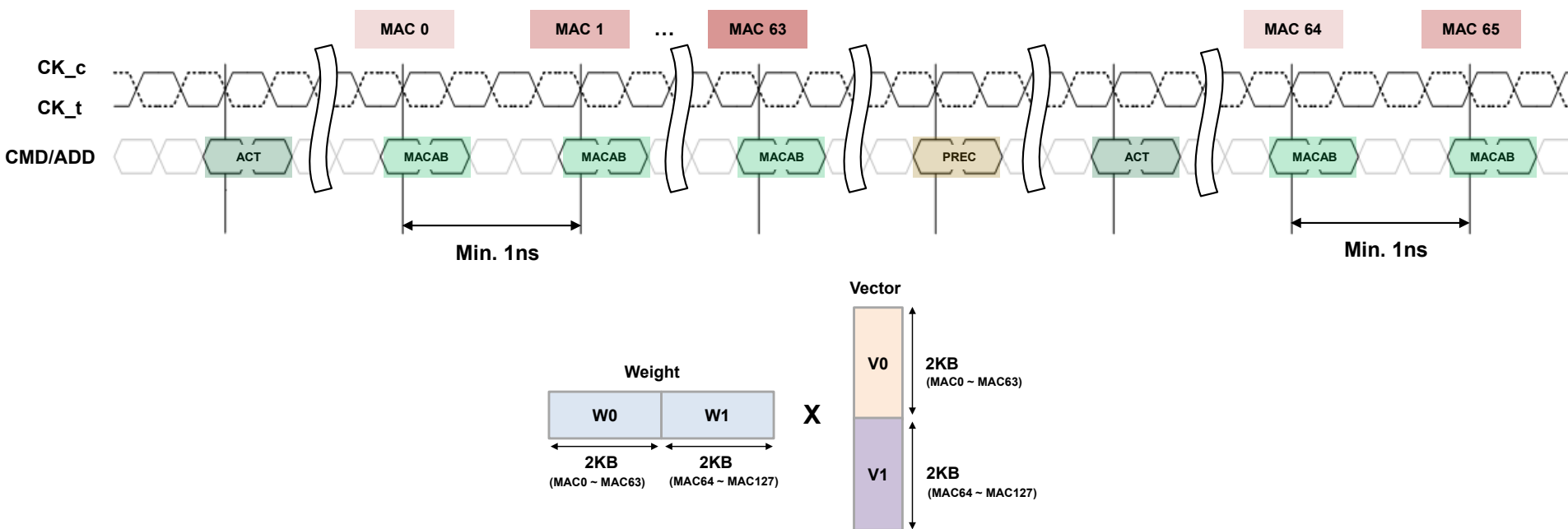
Type	CMD	Description
Bank Activation	ACT4, ACT16	Activate four/sixteen banks in parallel
	ACTAF4, ACTAF16	Activate rows storing activation function LUTs in four/sixteen banks in parallel
Compute	MACSB, MACAB, MACAB	Perform MAC in one/four/sixteen banks in parallel
	AF	Compute activation function in all banks
	EWML	Perform element-wise multiplication
Data	WRBK	Write to all activated banks in parallel
	WRGB	Write to Global Buffer
	RDMAC	Read from MAC result register
	RDAF	Read from activation function result register
	WRBIAS	Write bias to MAC result register
	RDCP	Copy data from a bank to Global Buffer
	WRCP	Copy data from Global Buffer to a bank



AiM: Accumulation Operation

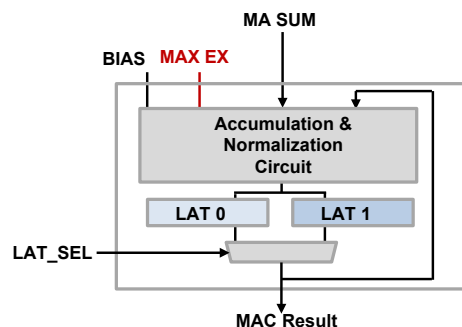
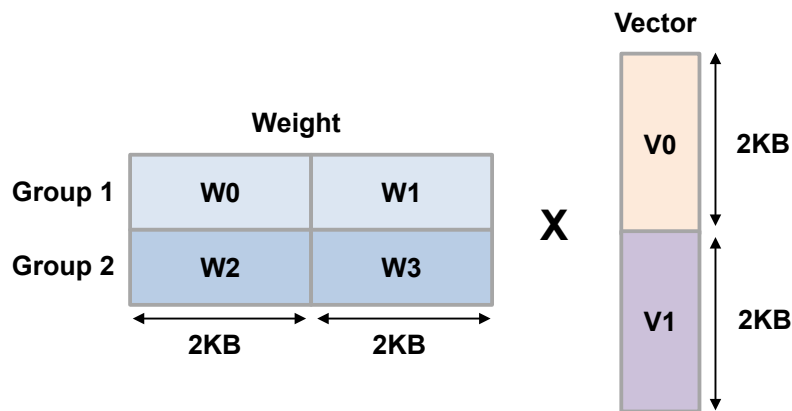
- Up to 64 MAC commands between ACT
 - 1 MAC every 1ns (t_{CCDS})
 - 64 x 256 bits = 2 KB (row size)

Type	CMD	Description
Bank Activation	ACT4, ACT16	Activate four/sixteen banks in parallel
	ACTAF4, ACTAF16	Activate rows storing activation function LUTs in four/sixteen banks in parallel
Compute	MACSB, MACAB, MACAB	Perform MAC in one/four/sixteen banks in parallel
	AF	Compute activation function in all banks
	EWML	Perform element-wise multiplication
Data	WRBK	Write to all activated banks in parallel
	WRGB	Write to Global Buffer
	RDMAC	Read from MAC result register
	RDAF	Read from activation function result register
	WRBIAS	Write bias to MAC result register
	RDCP	Copy data from a bank to Global Buffer
	WRCP	Copy data from Global Buffer to a bank



AiM: Multi-latch Operation

■ Save writes with selectable latches

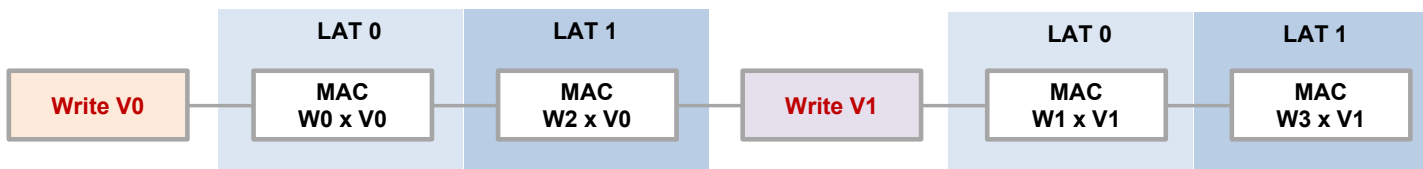


Accumulator with Two Selectable Latches Storing MAC Results

Type	CMD	Description
Bank Activation	ACT4, ACT16	Activate four/sixteen banks in parallel
	ACTAF4, ACTAF16	Activate rows storing activation function LUTs in four/sixteen banks in parallel
Compute	MACSB, MAC4B, MACAB	Perform MAC in one/four/sixteen banks in parallel
	AF	Compute activation function in all banks
Data	EWMUL	Perform element-wise multiplication
	WRBK	Write to all activated banks in parallel
	WRGB	Write to Global Buffer
	RDMAC	Read from MAC result register
	RDAF	Read from activation function result register
	WRBIAS	Write bias to MAC result register
	RDCP	Copy data from a bank to Global Buffer
	WRCP	Copy data from Global Buffer to a bank



Operation Sequence without Multi-Latch

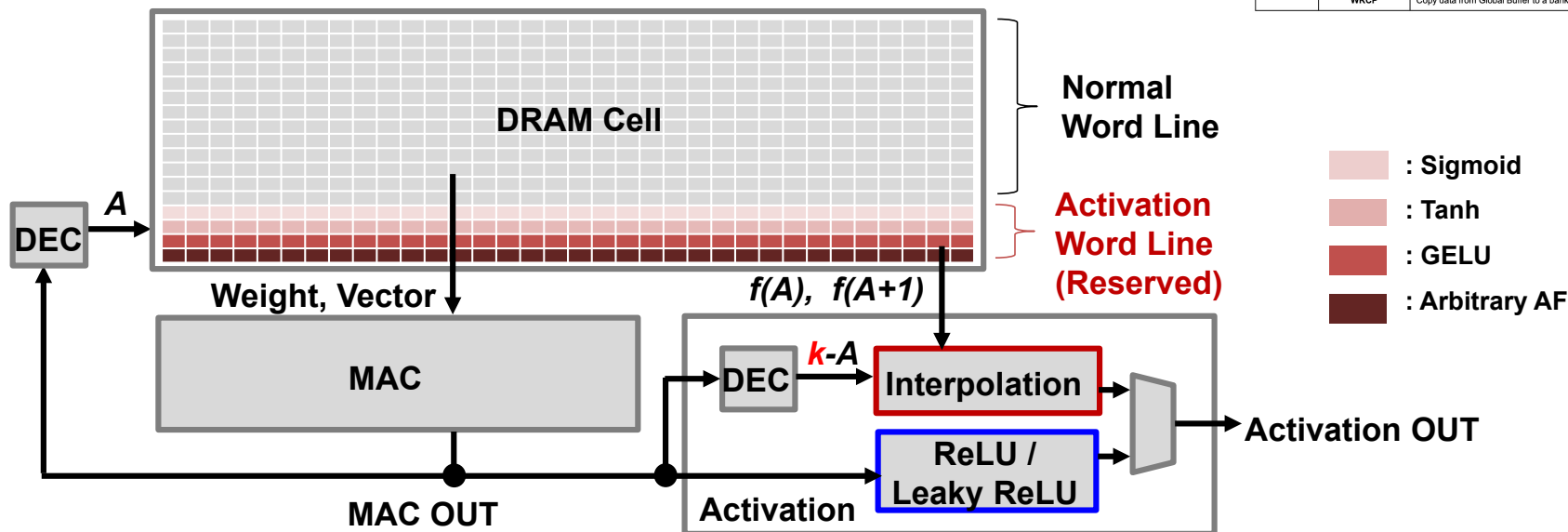


Operation Sequence with Multi-Latch

AiM: Activation Functions (I)

- Two types of activation functions
 - Calculation: ReLU, Leaky ReLU
 - With LUT: Sigmoid, GELU, Tanh, arbitrary AF

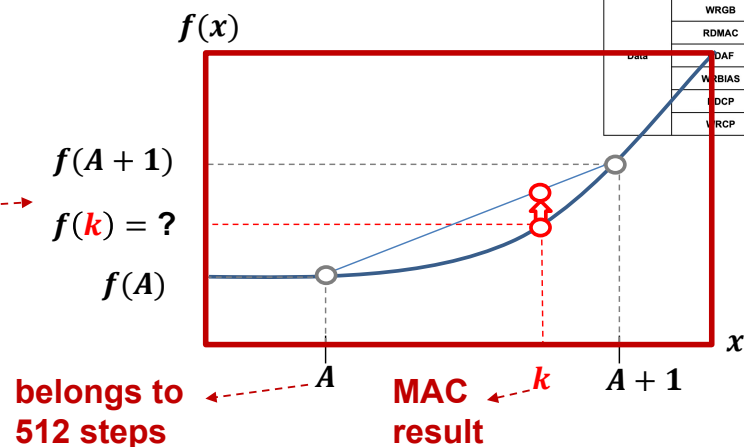
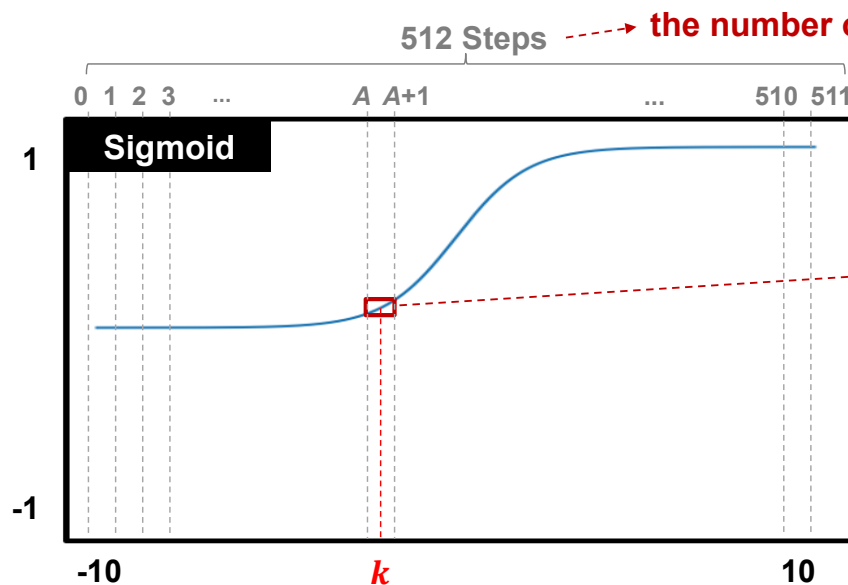
Type	CMD	Description
Bank Activation	ACT4, ACT16	Activate four/sixteen banks in parallel
	ACTAF4, ACTAF16	Activate rows storing activation function LUTs in four/sixteen banks in parallel
Compute	MACSB, MAC4B, MACAB	Perform MAC in one/four/sixteen banks in parallel
	AF	Compute activation function in all banks
	EWML	Perform element-wise multiplication
Data	WRBK	Write to all activated banks in parallel
	WRGB	Write to Global Buffer
	RDMAC	Read from MAC result register
	RDAF	Read from activation function result register
	WRBIAS	Write bias to MAC result register
	RDCP	Copy data from a bank to Global Buffer
	WRCP	Copy data from Global Buffer to a bank



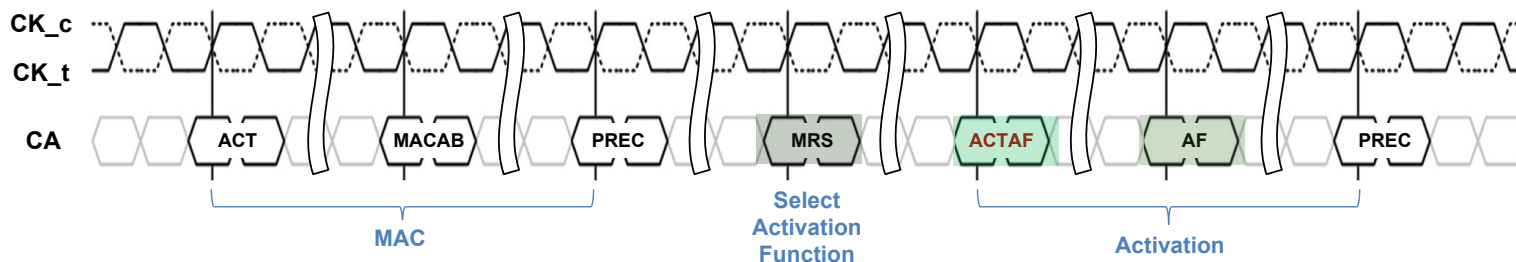
AiM: Activation Functions (II)

- LUT (512 values) + Linear interpolation
 - Mode set register (MSR) to select AF

Type	CMD	Description
Bank Activation	ACT4, ACT16	Activate four/sixteen banks in parallel
	ACTAF4, ACTAF16	Activate rows storing activation function LUTs in four/sixteen banks in parallel
Compute	MACSB, MAC4B, MACAB	Perform MAC in one/four/sixteen banks in parallel
	AF	Compute activation function in all banks
	EWML	Perform element-wise multiplication
	WRBK	Write to all activated banks in parallel
	WRGB	Write to Global Buffer
	RDMAC	Read from MAC result register
Data	DAF	Read from activation function result register
	WRBIAS	Write bias to MAC result register
	DCP	Copy data from a bank to Global Buffer
	VRCP	Copy data from Global Buffer to a bank



$$f(k) \doteq f(A) + \frac{k - A}{(A + 1) - A} \times (f(A + 1) - f(A))$$

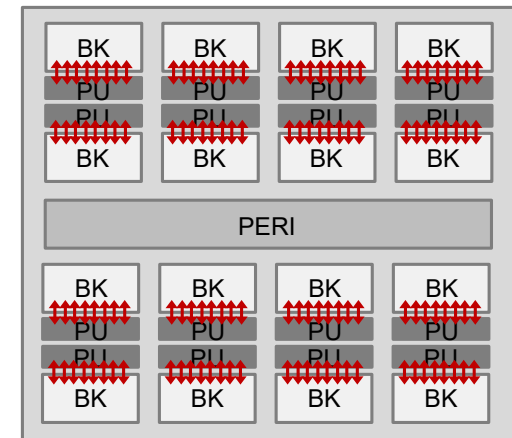


AiM: Key Feature Summary

■ Comparison table

	[1]	UPMEM PIM [2]	FIMDRAM [3, 7]	This work
DRAM Type	LPDDR4	DDR4	HBM2	GDDR6
Process	20nm	2x nm	20 nm	1y nm
Memory Density	8GB/chip (8H 8Gb mono die)	8GB/DIMM	6GB/cube (Buffer die + 4H 4Gb core-die with + 4H 8Gb core-die)	8Gb/chip (4Gb DDP)
Data Rate	3.2Gbps	2.4Gbps	2.4Gbps	16Gbps
Bandwidth	25.6GB/s per chip	19.2GB/s per DIMM	307GB/s per cube	64GB/s per chip
# of Processing Unit (PU)	2048 per chip (256 per die)	128 per DIMM (8 per chip)	128 per cube (32 per core-die)	32 per chip (16 per die)
Processing Operation Speed	250MHz	500MHz	300MHz	1GHz
1 PU Throughput	2 GOPS (250MHz x 8byte)	4 GOPS (500MHz x 8byte)	9.6 GFLOPS (300MHz x 32byte)	32 GFLOPS (1GHz x 32byte)
Total Throughput (1 PU Throughput x # of PU)	0.5 TOPS per chip (2 GOPS x 256)	0.5 TOPS per DIMM (4 GOPS x 128)	1.2 TFLOPS per cube (9.6 GFLOPS x 128)	1 TFLOPS per chip (32 GFLOPS x 32)
Operation precision	INT8	INT8	FP16	BF16
Supported Activation Functions	-	-	ReLU	Sigmoid, Tanh, GELU, ReLU, Leaky ReLU and Arbitrary AF

AIM, based on GDDR6,
exploiting bank parallelism



- [1] H. Shin, et al., IEEE TCADICS 2018,
[2] F. Devaux, IEEE Hot Chips Symp. 2019,
[3] Y.-C. Kwon et al., ISSCC 2021

Upcoming Lectures

- More real-world PIM architectures
- Programming PIM systems
- More on workload characterization for PIM suitability
- PUM architectures and prototypes

P&S Processing-in-Memory

Real-World Processing-in-Memory Architectures:
SK Hynix Accelerator-in-Memory

Dr. Juan Gómez Luna

Prof. Onur Mutlu

ETH Zürich

Fall 2022

15 November 2022